

**NCBI Molecular Biology Resources**

*Advances in Cell and Molecular Biology*  
American Society of Nephrology

Peter Cooper  
National Center for Biotechnology Information

NBRI

---

---

---

---

---

---

---

**Talk Outline**

- About NCBI
- Molecular Databases
  - GenBank
  - RefSeq
- Web Access
  - Entrez
  - BLAST
- Tour with MLH1
- Live Web Searches
  - Candidate gene by genetic markers
  - Human homologue of rat gene by BLAST

NBRI

---

---

---

---

---

---

---

**The National Center for Biotechnology Information (NCBI)**

- Created as a part of NLM in 1988
  - Establish public databases
  - Research in computational biology
  - Develop software tools for sequence analysis
  - Disseminate biomedical information
- Tools: BLAST(1990), Entrez (1992)
- GenBank (1992)
- Free MEDLINE (PubMed, 1997)
- Human genome (2001)

NBRI

---

---

---

---

---

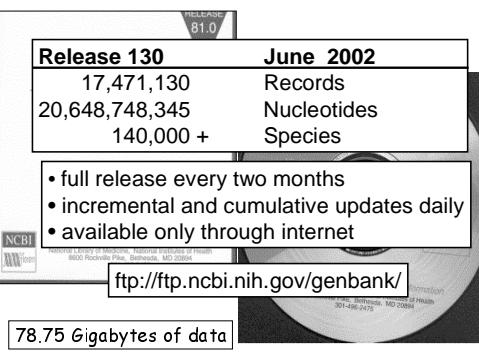
---

---

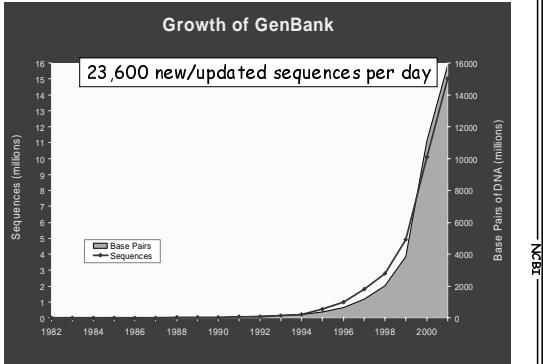
## Molecular Databases

- Primary Databases
  - Original submissions by experimentalists
  - Database staff organize but don't add additional information
    - Example: GenBank
- Derivative Databases
  - Human curated
    - compilation and correction of data
    - Example: SWISS-PROT, NCBI RefSeq mRNA
  - Computationally Derived
    - Example: UniGene
  - Combinations
    - Example: NCBI Genome Assembly

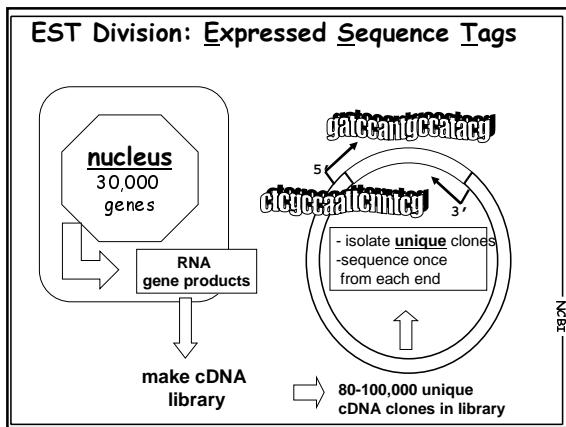
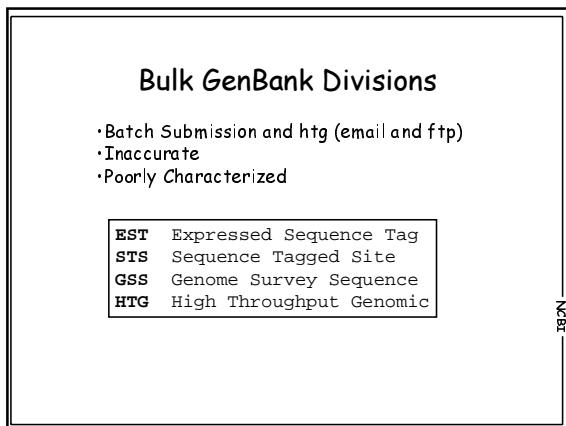
## GenBank: NCBI's Primary Sequence Database

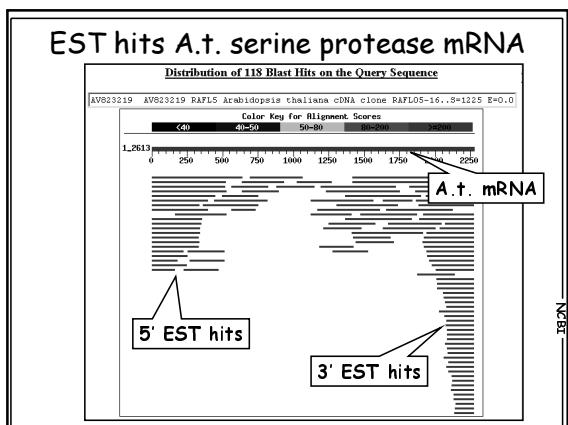
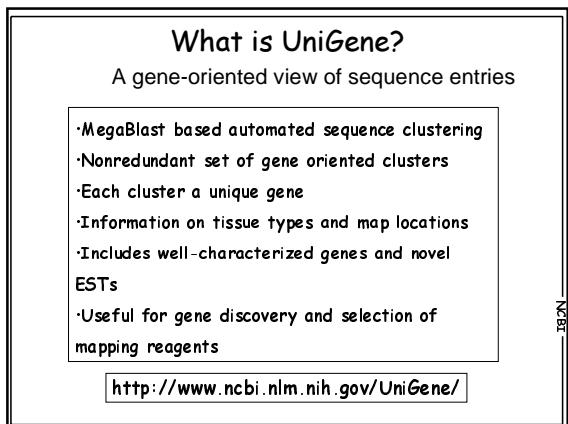
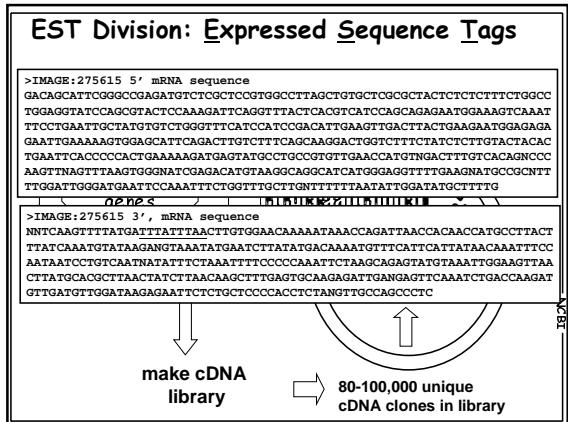


## Growth of GenBank



| GenBank Divisions       |                                     |     |     |     |     |
|-------------------------|-------------------------------------|-----|-----|-----|-----|
| Bulk Sequence Divisions |                                     |     |     |     |     |
| PAT                     | Patent                              |     |     |     |     |
| EST                     | Expressed Sequence Tags (142 files) |     |     |     |     |
| STS                     | Sequence Tagged Sites               |     |     |     |     |
| GSS                     | Genome Survey Sequences (48 files)  |     |     |     |     |
| HTG                     | High Throughput Genome (26 files)   |     |     |     |     |
| HTC                     | High Throughput cDNA                |     |     |     |     |
| CON                     | Contig                              |     |     |     |     |
| Traditional Divisions   |                                     |     |     |     |     |
| BCT                     | INV                                 | MAM | PHG | PLN | PRI |
| ROD                     | SYN                                 | UNA | VRL | VRT |     |

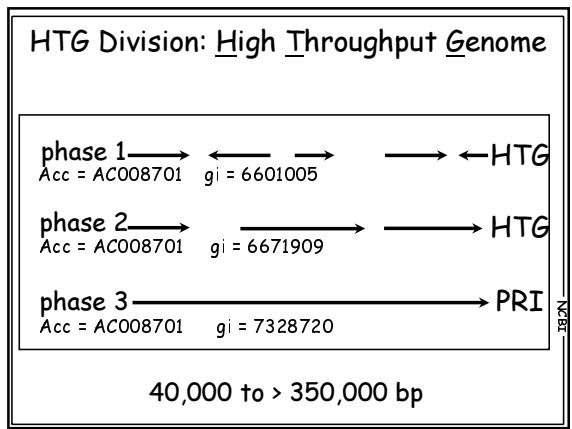
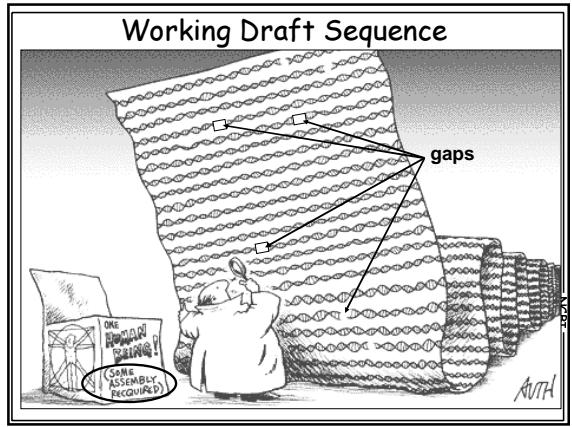
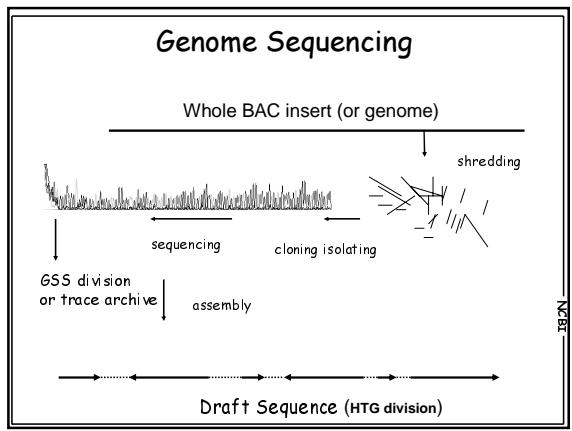




| Arabidopsis UniGene Statistics      |                                      |
|-------------------------------------|--------------------------------------|
| 39,855                              | mRNAs + gene CDSs                    |
| 87,006                              | EST, 3'reads                         |
| 42,137                              | EST, 5'reads                         |
| + 32,571                            | EST, other/unknown                   |
| <hr/>                               |                                      |
| 201,569                             | total sequences in clusters          |
| <br>Final Number of Clusters (sets) |                                      |
| <hr/>                               |                                      |
| 26,808                              | sets total 115,000,000 bp            |
| 25,474                              | sets contain at least one known gene |
| 17,654                              | sets contain at least one EST        |
| 16,326                              | sets contain both genes and ESTs     |

| Hs UniGene Statistics               |                                      |
|-------------------------------------|--------------------------------------|
| 73,419                              | mRNAs + gene CDSs                    |
| 1,181,855                           | EST, 3'reads                         |
| 1,461,928                           | EST, 5'reads                         |
| + 616,609                           | EST, other/unknown                   |
| <hr/>                               |                                      |
| 3,333,811                           | total sequences in clusters          |
| <br>Final Number of Clusters (sets) |                                      |
| <hr/>                               |                                      |
| 98,816                              | sets total 3,000,000 base pairs      |
| 22,431                              | sets contain at least one known gene |
| 97,618                              | sets contain at least one EST        |
| 21,233                              | sets contain both genes and ESTs     |

| UniGene Collections Apr. 2002 |                            |
|-------------------------------|----------------------------|
| Animals                       | Sequences Clusters         |
| <i>Homo sapiens</i>           | human 3,333,811 98,816     |
| <i>Mus musculus</i>           | mouse 2,274,640 86,897     |
| <i>Rattus norvegicus</i>      | rat 308,877 59,882         |
| <i>Danio rerio</i>            | zebrafish 159,261 14,893   |
| <i>Bos taurus</i>             | cow 122,503 9,303          |
| <i>Xenopus laevis</i>         | frog 120,489 16,489        |
| <i>Anopheles gambiae</i>      | mosquito 42,590 2,414      |
| Plants                        |                            |
| <i>Arabidopsis thaliana</i>   | thale cress 202,099 26,794 |
| <i>Oryza sativa</i>           | rice 77,376 15,283         |
| <i>Triticum aestivum</i>      | wheat 35,387 3,091         |
| <i>Hordeum vulgare</i>        | barley 108,658 6,984       |
| <i>Zea mays</i>               | maize (corn) 108,030 9,889 |



## RefSeq: NCBI's Derivative Sequence Database

- Curated transcripts and proteins
  - reviewed
  - human, mouse, rat, fruit fly, zebrafish, arabidopsis
- Human model transcripts and proteins
- Assembled Genomic Regions (contigs)
  - draft human genome
  - mouse genome
- Chromosome records
  - microbial
  - organelle

## RefSeq Resource

- Genome Oriented Resource
  - A sequence for each macromolecule Central Dogma: Chromosome, mRNA, preprotein, mature protein
  - Linked on a residue by residue basis
  - Objectively non-redundant and comprehensive
- Curated Resource
  - Authoritative source by genome
  - Derived from GenBank but corrected, merged, extended
  - Publicly distributed, Entrez Genomes Web site
- Reagents for Genome Annotation and Analysis
- Substrate for Functional Genomics

## The RefSeq Accession Numbers

### NCBI Reference Sequences

#### mRNAs and Proteins

NM\_123456 Curated mRNA  
NP\_123456 Curated Protein  
XM\_123456 Predicted Transcript (human)  
XP\_123456 Predicted Protein (human)

human  
mouse  
rat  
fruit fly  
zebrafish  
Arabidopsis

#### Gene Records

NG\_123456 Reference Genomic Sequence (human)

#### Assemblies

NT\_123456 Contig (Mouse and Human)  
NW\_123456 Super Contig (Mouse)  
NC\_123455 Chromosome (Microbial, Arabidopsis )

**GenBank Sequences: human CFTR**

Show 100 ▾ Items 1-83 of 83 One page

1: M86631 Related Sequences, OMIM, Protein, PubMed, Taxonomy  
Home sapiens (clone ST-13-5(9/16)) cystic fibrosis transmembrane conductance regulator (CFTR) gene; 3' end intron 17B; complete exon 18; complete intron 18  
gi|180296|gb|M86631.1|HUMCFTRI|180296|

2: S64699 Related Sequences, OMIM, Protein, PubMed, Taxonomy  
Home sapiens cystic fibrosis transmembrane conductance regulator isoform 36 (CFTR) mRNA, partial cds  
gi|408283|gb|S64699.1|S64699|408285|

3: AL121762 Related Sequences, Taxonomy  
Human DNA sequence from clone RP4-6|10C12 on chromosome 20 Contains the 2' end of a novel gene, a putative novel gene, a pseudogene similar to part of the cystic fibrosis transmembrane conductance regulator (CFTR), four CpG islands, ESTs, STSs and GSSs, complete sequence  
gi|8574423|emb|AL121762.1|HSM610C12|8574423|

4: AH006034 Clin., Protein, PubMed, Taxonomy  
Human cystic fibrosis transmembrane conductance regulator (CFTR) gene  
gi|306537|gb|AH006034.1|SEG\_HUMCFTRA|306537|

5: M55131 Related Sequences, ProbeSet, OMIM, Protein, PubMed, Taxonomy  
Human cystic fibrosis transmembrane conductance regulator (CFTR) gene, exon 24  
gi|306536|gb|M55131.1|HUMCFTRAA24|306536|

6: M55130 Related Sequences, Protein, PubMed, Taxonomy  
Human cystic fibrosis transmembrane conductance regulator (CFTR) gene, exon 23

**Curated RefSeq Records: NM\_, NP\_**

LOCUS NM\_000492 6159 bp mRNA PRI 26-JUL-1999  
DEFINITION Homo sapiens cystic fibrosis transmembrane conductance regulator(CFTR) mRNA.  
ACCESSION NM\_000492 RefSeq Nucleotide

LOCUS NP\_000483 1480 aa PRI 26-JUL-1999  
DEFINITION cystic fibrosis transmembrane conductance regulator.  
ACCESSION NP\_000483 RefSeq Protein  
PID g4502785  
VERSION NP\_000483.1 GI:4502785  
DBSOURCE RefSeq: accession NM\_000492.1

COMMENT REFSEQ: This reference sequence was derived from M55131.  
PROVISIONAL RefSeq: This is a provisional reference sequence record that has not yet been subject to human review. The final curated reference sequence record may be somewhat different from this one.

**Curated RefSeq Records: NM\_, NP\_**

LOCUS NM\_000492 6159 bp mRNA PRI 26-JUL-1999  
DEFINITION Homo sapiens cystic fibrosis transmembrane conductance regulator(CFTR) mRNA  
ACCESSION NM\_000492 RefSeq Nucleotide  
REFSEQ: This reference sequence was derived from M28668.1, M55131.1. e  
On Feb 17, 2000 this sequence version replaced gi|4502784. 1999  
Summary: Cystic fibrosis transmembrane conductance regulator is member 7 of the ATP-binding cassette sub-family C. The protein functions as a chloride channel and controls the regulation of other transport pathways. Mutations in this gene cause the autosomal recessive disorder, cystic fibrosis (CF) and congenital bilateral aplasia of the vas deferens (CBAD). Alternative splice variants have been described, many of which result from mutations in the CFTR gene.  
COMPLETENESS: full length. Reviewed n

COMMENT REFSEQ: This reference sequence was derived from M55131.  
PROVISIONAL RefSeq: This is a provisional reference sequence record that has not yet been subject to human review. The final curated reference sequence record may be somewhat different from this one.

**Alignment Generated Transcripts: XM\_, XP\_**

Genomic: [gi|13631467|ref|NT\\_007935.3|Hs7\\_8092](#) Homo sapiens chromosome 7 working draft sequence segment

mRNA: [gi|6995995|ref|NM\\_000492.2|](#) Homo sapiens cystic fibrosis transmembrane conductance regulator, ATP-binding cassette (sub-family C, member 7) (CFTR), mRNA

Alignment is on plus strand of genomic sequence and on plus strand of mRNA sequence  
mRNA coverage: 100%  
Overall percent identity: 99.9%

634416 |-----|-----|-----|-----|-----||823113

---

---

---

---

---

---

---

---

---

---

**Alignment Generated Transcripts: XM\_, XP\_**

Genomic: [gi|13631467|ref|NT\\_007935.3|Hs7\\_8092](#) 7 working draft sequence

mRNA: [gi|6995995|ref|NM\\_000492.2|](#) transmembrane conductance regulator, ATP-binding cassette (sub-family C, member 7) (CFTR), mRNA

Alignment is on plus strand of genomic sequence  
mRNA coverage: 100%  
Overall percent identity: 99.9%

634416 |-----|-----|-----|-----|-----||823113

| Exon    | Genomic coordinates | mRNA coordinates | length | identity | mismatches | gaps | Donor site | Acc. site |
|---------|---------------------|------------------|--------|----------|------------|------|------------|-----------|
| Exon 1  | 634415-634599       | 1-185            | 185    | 100.0%   | 0          | 0    | d          | a         |
| Exon 2  | 658705-658815       | 186-296          | 111    | 100.0%   | 0          | 0    | d          | a         |
| Exon 3  | 663486-663594       | 297-405          | 109    | 100.0%   | 0          | 0    | d          | a         |
| Exon 4  | 68351-68556         | 406-621          | 216    | 100.0%   | 0          | 0    | d          | a         |
| Exon 5  | 688728-688817       | 622-711          | 90     | 100.0%   | 0          | 0    | d          | a         |
| Exon 6  | 689700-689863       | 712-875          | 164    | 100.0%   | 0          | 0    | d          | a         |
| Exon 7  | 691000-691125       | 876-1001         | 126    | 100.0%   | 0          | 0    | d          | a         |
| Exon 8  | 694552-694798       | 1002-1248        | 247    | 100.0%   | 0          | 0    | d          | a         |
| Exon 9  | 696468-696560       | 1249-1341        | 93     | 100.0%   | 0          | 0    | d          | a         |
| Exon 10 | 703093-703275       | 1342-1524        | 183    | 100.0%   | 0          | 0    | d          | a         |
| Exon 11 | 713916-714107       | 1525-1716        | 192    | 99.5%    | 1          | 0    | d          | a         |
| Exon 12 | 742191-742285       | 1717-1811        | 95     | 100.0%   | 0          | 0    | d          | a         |
| Exon 13 | 744805-744891       | 1812-1898        | 87     | 100.0%   | 0          | 0    | d          | a         |

---

---

---

---

---

---

---

---

---

---

**Alignment Generated Transcripts: XM\_, XP\_**

Genomic: [gi|13631467|ref|NT\\_007935.3|Hs7\\_8092](#) 7 working draft sequence

mRNA: [gi|6995995|ref|NM\\_000492.2|](#) transmembrane conductance regulator, ATP-binding cassette (sub-family C, member 7) (CFTR), mRNA

Alignment is on plus strand of genomic sequence  
mRNA coverage: 100%  
Overall percent identity: 99.9%

634416 |-----|-----|-----|-----|-----||823113

Exon 11: 713916-714107 (genomic); 1525-1716 (mRNA)

The alignment shows the genomic sequence (top) and mRNA sequence (bottom). A mismatch is highlighted at position 1525 of the mRNA.

```

TTATTTCCAGACTTCACCTTCTAAATGCGATTATGGAGAGACTGGAGCCTTCAAGGGGTAAT
| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
ACTTCACCTCTAATGTGATTATGGAGAGACTGGAGCCTTCAGAGGGTTAATGCCTGGAC
T S L L M H I M G E L E F S E G K
| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
AATTAAGCAGCATGCAGTGGAAGGATAAGTTTCTGGATAATGCCCTGGAGATATGCCTGGAC
I K H S G R I S F C S Q F S W I M F G T
| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
AATTAAGCAGCATGCAGTGGAAGGAAATTTCATTCTGTCTCAGTTTCTGGATATATGCCTGGAC
I K E N I I F G V S Y D E Y R Y R S V I
| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
CATTAAAGAAAAATATCATCTTTGGTTCTATGTGATGATATAGATACAGAACGGTCAT
| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
CATTAAAGAAAAATATCATCTTTGGTTCTATGTGATGATATAGATACAGAACGGTCAT
I K E N I I F G V S Y D E Y R Y R S V I
| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
CAAAGCATGCCAAGTAGAACAGGTAAAGAAACT
| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
CAAAGCATGCCAAGTAGAACAGGTAAAGAAACT
K A C Q L E E
| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |

```

---

---

---

---

---

---

---

---

---

---

**Alignment Generated Transcripts: XM\_, XP\_**

Exon 11: 713916-714107 (genomic); 1525-1716 (mRNA)

```

  TTATTTCCAGACTTCACCTTCTAATGGGATTATGGGAGAACTGGAGCCTTCAGAGGGTAA
  ||||| ||||| ||||| ||||| ||||| ||||| ||||| ||||| ||||| ||||| ||||| |||||
  ACTTCACCTTCTAATGGGATTATGGGAGAACTGGAGCCTTCAGAGGGTAA
  T S L H I M G E L E F S E G K
  AATTAAGCAGTGGGAGATACTTCAGTTCTGTCAGTTCTGGATTATGCGCTGGCAC
  ||||| ||||| ||||| ||||| ||||| ||||| ||||| ||||| ||||| ||||| |||||
  AATTAAGCAGTGGGAGATACTTCAGTTCTGTCAGTTCTGGATTATGCGCTGGCAC

```

mismatch

LOCUS XM\_004980 6128 bp mRNA PRI 16-NOV-2000  
DEFINITION Homo sapiens cystic fibrosis transmembrane conductance regulator, ATP-binding cassette (sub-family C, member 7) (CFTR), mRNA.  
ACCESSION XM\_004980  
VERSION XM\_004980.3 GI:13631444

|                    |                                  |
|--------------------|----------------------------------|
| reference sequence | Exon 11                          |
| mRNA c             | CAAAGCATGCCAAGTAGAAAGGGTAAGAAACT |
| Overall p          |                                  |
| 634416             | Exon 12                          |
|                    | CAAAGCATGCCAAGTAGAAAGGG          |
|                    | K A C Q L E E                    |
|                    | Exon 13                          |

**RefSeq Contig: NT\_**

LOCUS NT\_007933 22893611 bp DNA linear CON 13-MAY-2002  
DEFINITION Homo sapiens chromosome 7 working draft sequence segment.  
ACCESSION NT\_007933  
VERSION NT\_007933.9 GI:20543589  
KEYWORDS .  
SOURCE .  
ORGANISM Homo sapiens  
Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;  
Mammalia; Eutheria; Primates; Catarrhini; Hominidae; Homo.  
REFERENCE 1 (bases 1 to 22893611)  
AUTHORS NCBI Annotation Project.  
TITLE Direct Submission  
JOURNAL Submitted (09-MAY-2002) National Center for Biotechnology  
Information, NIH, Bethesda, MD 20894, USA  
COMMENT GENOME ANNOTATION REFSEQ: NCBI contigs are derived from assembled  
genomic sequence data. They may include both draft and finished  
sequence.  
On May 13, 2002 this sequence version replaced gi:18568843.  
COMPLETENESS: not full length.

**RefSeq Contig: NT\_**

LOCUS NT\_007933 22893611 bp DNA linear CON 13-MAY-2002  
DE  
AC gene 9640781..9829479  
VE /gene="CFTR"  
KE /note="CF; MRP7; ABC35; ABCC7"  
SC /db\_xref="LocusID:1080"  
RE /db\_xref="MM:602421"  
MRNA join(9640781..9640965,9665071..9665181,9669852..9669960,  
9691717..9691932,9695094..9695183,9696066..9696229,  
9697366..9697401,9700491..9702834..9702926,  
9703044..9703141..9703245..9703351..9703451,  
9751171..9751257,9752752..9753475,9755748..9755876,  
9763644..9763681,9764350..9764600,9767492..9767571,  
9771327..9771487,9772399..9772626,9775421..9775531,  
9788340..9788588,9803256..9803411,9813660..9813749,  
9825506..9825678,9826277..9826382,9827726..9829479)  
/genes="CFTR"  
/product="cystic fibrosis transmembrane conductance  
regulator, ATP-binding cassette (sub-family C, member 7)"  
/note="Derived by automated computational analysis using  
gene prediction method: BLAST. Supporting evidence  
includes similarity to: 2 mRNAs"  
/transcript\_id="XM\_004980.4"  
/db\_xref="GI:14753226"  
/db\_xref="LocusID:1080"  
/db\_xref="MM:602421"

**RefSeq Contig: NT**

|   |   |             |     |        |                 |
|---|---|-------------|-----|--------|-----------------|
| LOCUS   | NT_007933   | 22893611 bp | DNA | linear | CON 13-MAY-2002 |
| DEFINITION  | Mus musculus WIFeb01_97.  |             |     |        |                 |
| ACCESSION   | NW_000272   |             |     |        |                 |
| VERSION   | NW_000272.1 GI:20915425   |             |     |        |                 |
| KEYWORDS  | .   |             |     |        |                 |
| SOURCE  | house mouse.  |             |     |        |                 |
| ORGANISM  | Mus musculus  |             |     |        |                 |
| Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi; Mammalia; Eutheria; Rodentia; Sciurognathia; Muridae; Murinae; Mus. |   |             |     |        |                 |
| REFERENCE   | 1 (bases 1 to 27755676)   |             |     |        |                 |
| AUTHORS   | NCBI Annotation Project.  |             |     |        |                 |
| TITLE   | Direct Submission   |             |     |        |                 |
| JOURNAL   | Submitted (13-MAY-2002) National Center for Biotechnology Information, NIH, Bethesda, MD 20894, USA   |             |     |        |                 |
| COMMENT   | GENOME ANNOTATION REFSEQ: NCBI contigs are derived from assembled genomic sequence data. They may include both draft and finished sequence. |             |     |        |                 |
| /db_xref="LocusID:1080"   |   |             |     |        |                 |
| /db_xref="MIM:602421"   |   |             |     |        |                 |

**RefSeq WGS Supercontig: NW**

|   |   |             |     |        |                 |
|---|---|-------------|-----|--------|-----------------|
| LOCUS   | NW_000272   | 27755676 bp | DNA | linear | CON 04-JUN-2002 |
| DEFINITION  | Mus musculus WGS supercontig Mm6_WIFeb01_97.  |             |     |        |                 |
| ACCESSION   | NW_000272   |             |     |        |                 |
| VERSION   | NW_000272.1 GI:20915425   |             |     |        |                 |
| KEYWORDS  | .   |             |     |        |                 |
| SOURCE  | house mouse.  |             |     |        |                 |
| ORGANISM  | Mus musculus  |             |     |        |                 |
| Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi; Mammalia; Eutheria; Rodentia; Sciurognathia; Muridae; Murinae; Mus. |   |             |     |        |                 |
| REFERENCE   | 1 (bases 1 to 27755676)   |             |     |        |                 |
| AUTHORS   | NCBI Annotation Project.  |             |     |        |                 |
| TITLE   | Direct Submission   |             |     |        |                 |
| JOURNAL   | Submitted (15-MAY-2002) National Center for Biotechnology Information, NIH, Bethesda, MD 20894, USA   |             |     |        |                 |
| COMMENT   | GENOME ANNOTATION REFSEQ: NCBI contigs are derived from assembled genomic sequence data. They may include both draft and finished sequence. |             |     |        |                 |

**RefSeq WGS Supercontig: NW**

|          |  |             |     |        |                 |
|----------|--|-------------|-----|--------|-----------------|
| LOCUS    | NW_000272  | 27755676 bp | DNA | linear | CON 04-JUN-2002 |
| FEATURES | Location/Qualifiers  |             |     |        |                 |
| source   | 1..27755676 /organism="Mus musculus" /strain="C57BL/6J" /db_xref="taxon:10090" /chromosome="6"   |             |     |        |                 |
| gene     | 6728..13494 /genes="LOC243722" /db_xref="InterPro:243722"  |             |     |        |                 |
| mRNA     | join(6728..6872,7967..8218,11029..11098,11202..11335,11368..11503,11641..11766,12089..12225,12970..13494) /genes="LOC243722" /product="similar to nuclear matrix transcription factor" /note="Derived by automated computational analysis using gene prediction method: GenomicScan." /transcript_id="XM_145115.1" /db_xref="GI:20915205" /db_xref="InterPro:243722" |             |     |        |                 |
| CDS      | join(6786..6972,7967..8218,11029..11098,11202..11335,11368..11503,11641..11766,12089..12225,12970..13494) /genes="LOC243722" /codon_start=1 /protein_id="XP_145115.1" /db_xref="GI:20915206" /db_xref="InterPro:243722"  |             |     |        |                 |

| RefSeq WGS Supercontig: NW_  |           |             |     |        |                 |
|--|-----------|-------------|-----|--------|-----------------|
| LOCUS  | NW_000272 | 27755676 bp | DNA | linear | CON 04-JUN-2002 |
| <b>FEATURES Location/Qualifiers</b>  |           |             |     |        |                 |
| source 1..27755676 /organism="Mus musculus"<br>/strain="C57BL/6J"<br>/db_xref="taxon:10090"<br>/chromosome="6"<br>gene 6728..13494   |           |             |     |        |                 |
| <b>CONTIG</b>  |           |             |     |        |                 |
| join(CAA01202740..1..4651,gap(100),CAA01125491..1..11772,<br>gap(189),CAA01120088..1..4194,gap(122),CAA01145531..1..18939,<br>gap(243),CAA01174720..1..20193,gap(100),CAA01148186..1..4432,<br>gap(100),CAA01138802..1..26173,gap(100),CAA01149691..1..10320,<br>gap(100),CAA01147818..1..1458,gap(1004),CAA01147815..1..6125,<br>gap(749),CAA01147812..1..1561,gap(100),CAA01074766..1..2776,<br>gap(894),CAA01074767..1..8180,gap(183),CAA01138806..1..2425,<br>gap(100),CAA01144336..1..11048,gap(100),CAA01191698..1..2516,<br>gap(330),CAA01048494..1..6532,gap(100),CAA01048497..1..1666,<br>gap(1506),CAA01044252..1..76338,gap(100),CAA01171120..1..13327,<br>gap(100),CAA01044255..1..5744,gap(100),CAA01222036..1..8808,<br>gap(100),CAA01159465..1..4183,gap(271),CAA01044258..1..20222,<br>gap(504),CAA01202743..1..2875,gap(87),CAA01114772..1..2380,<br>AC00688..1..6379..AC00790..1..17369<br>CAA01044276..1..58306..74574..gap(100),CAA01044280..1..10478,<br>gap(100),CAA01044282..1..1857..gap(1569),CAA01044288..1..126746, |           |             |     |        |                 |

## RefSeq: NCBI's Derivative Sequence Database

- Curated transcripts and proteins
  - reviewed
  - human, mouse, rat, fruit fly, zebrafish, arabidopsis
- Human model transcripts and proteins
- Assembled Genomic Regions (contigs)
  - draft human genome
  - mouse genome
- Chromosome records
  - microbial
  - organelle

## Integrated WWW Access: BLAST and Entrez

NCBI National Center for Biotechnology Information

Search Nucleotide [for]

What does NCBI do?

Established in 1988 as a national resource for molecular biology information, NCBI creates public databases, conducts research in computational biology, develops software tools for analyzing genome data, and disseminates biomedical information - all for the better understanding of human health and disease. [More](#)

Mouse Genome

Resources: explore tools for manipulating the mouse genome.

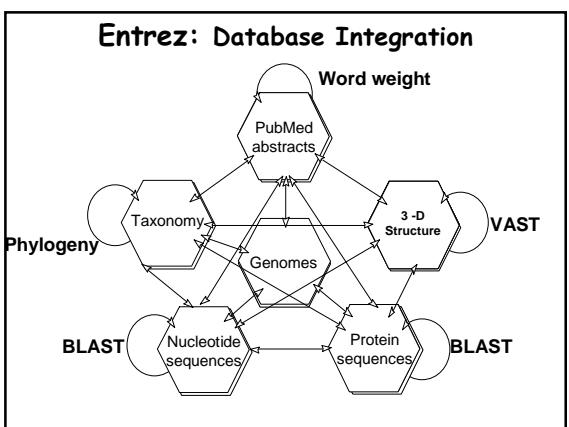
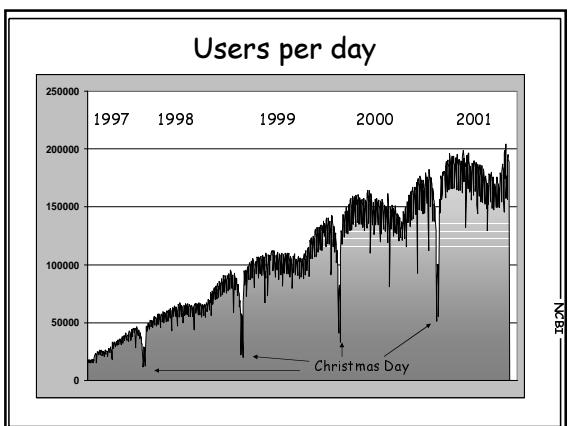
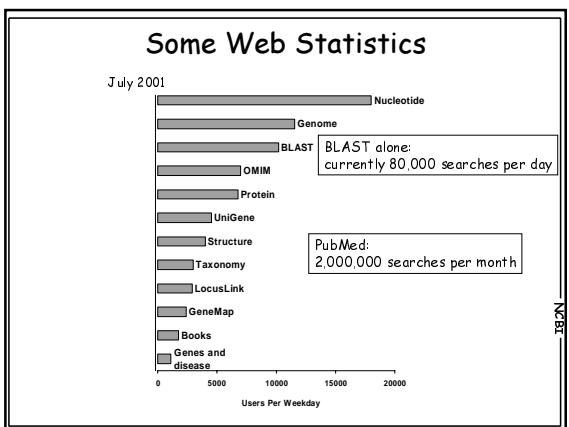
Try these: [Find New](#) [Sequence](#) [Human Mouse](#) [Homology](#) [Map](#)

Draft Human Genome

Explore human genome resources or browse the human genome as it was assembled by the ENSEMBL project.

Genes, genomes, and phylogenetic patterns

A recent whole genome analysis of the human genome indicates that Methoprotococcus Mendeleyevii is an archaeal phylogenetic outlier with other methanogens and is not a deep branching species close to the root of the tree. These results emphasize the importance of genome sequencing for accurate reconstruction of phylogeny. [View](#)



## The Draft Human Genome

NB#

---

---

---

---

---

---

---

---

---

---

---

## Human Genome Resources Find a gene by ...

NB#

---

---

---

---

---

---

---

---

---

---

---

## Other Genomes

NB#

---

---

---

---

---

---

---

---

---

---

---

**Other Genomes**

**scale science** [Organism-specific resources:](#)

- Fruit fly
- Human
- Human Genome Project
- Malaria parasite
- Microbial Genomes
- Plant Genomes Central
- Retroviruses
- Zebrafish

**The SNP Database**

Singlenucleotide polymorphisms (SNPs) are the most common genetic variations and occur once every 100 to 300 bases. It is expected

**Meet the Team**

**Other Genomes**

**Strain Resources**

**Microbial genomes**

**Human**

**Malaria parasite**

**Mouse**

**Plant Genomes Central**

**Retroviruses**

**Other Genomes**

**Strain Resources**

**Microbial genomes**

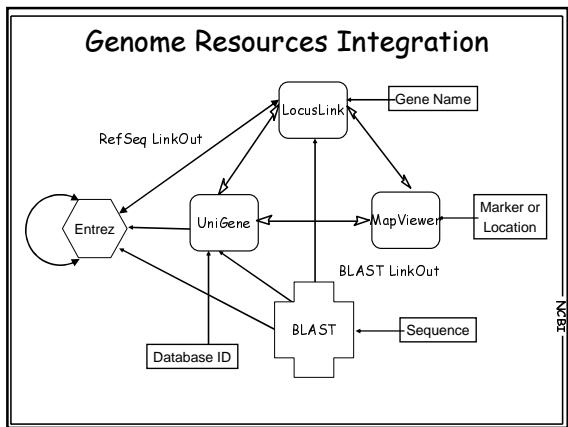
**Human**

**Malaria parasite**

**Mouse**

**Plant Genomes Central**

**Retroviruses**



**WWW Entrez**

**Human DNA Mismatch Repair Protein**



**NCBI SNP**

ENTREZ SNP

Search... Filter... Limits Previous/Next History Clipboard Details

Display: Fresh Default Sort Save Tool Clip Add One page

Entrez SNP (bulk)

SNP ID: 633129922

Nucleotide Protein

LocusLink via analysis of contig annotation: MLH11 mutL homolog 1, colon cancer, nonpolyposis type 2 (E. coli)

Gene Model (contig mRNA transcript) information from genome sequence for NM\_000249

Legend

| Contig accession | Contig position | Protein accession | Function         | dbSNP | Protein Codon         | Amino acid residue position | Position |   |    |
|------------------|-----------------|-------------------|------------------|-------|-----------------------|-----------------------------|----------|---|----|
| NM_000249        | 2631063         | NP_000249         | contig reference | A     | Ile [I]               | 1                           | 32       |   |    |
|                  |                 |                   |                  |       | non synonymous change | G                           | Val [V]  | 1 | 32 |

9: nt|NM\_000249| LocusLink, Nucleotide, Protein, Published HGVS

10: nt|NM\_000249| LocusLink, Nucleotide, Protein, Published HGVS

11: nt|NM\_000249| LocusLink, Nucleotide, Protein, Published HGVS

Query: g455752 (NM\_000249) mutL homolog 1; mutL (E. coli) homolog 1; mutL (E. coli) homolog 1 (colon cancer, nonpolyposis type 2) [Homo sapiens]

Matching #453989 631299 730028 741681 1079787 13805126

Allied Common Tree Taxonomy Report 3D structures CDD-Search GBL

6 BLAST hits to 87 unique species. Sort by taxonomy proximity.

Archaea Bacteria Metazoa Fungi Plants Viruses Other Eukaryotes

Keep only Cut-Off 100 Select Reset

**BlastLink**

746 aa

| SCORE | P   | ACCESSION | GI       | N  | ORGANISM  |
|-------|-----|-----------|----------|----|---|
| 3565  | 0.6 | AAA17774  | 4165462  | 35 | <i>Homo sapiens</i>                               |
| 3358  | 2.0 | AAK38906  | 1524119  | 13 | <i>Rattus norvegicus</i>                          |
| 3358  | 2.0 | AAK38906  | 1524119  | 13 | <i>Rattus norvegicus</i>                          |
| 1687  | 0   | AAA00335  | 3126781  | -  | <i>Anopheles gambiae str. PEST</i>                |
| 265   | 0   | AAK37993  | 1340773  | -  | <i>Arabidopsis thaliana</i>                       |
| 3143  | 3   | AAK37993  | 1340773  | -  | <i>Arabidopsis thaliana</i>                       |
| 3262  | 4   | CAAB9903  | 825572   | 11 | <i>Seecarophyton cerevisiae</i>                   |
| 3262  | 4   | CAAB9903  | 825572   | 11 | <i>Seecarophyton cerevisiae</i>                   |
| 1014  | 7   | CAB17283  | 3800333  | -  | <i>Caenorhabditis elegans</i>                     |
| 228   | 0   | AAK74000  | 1337782  | -  | <i>Oryza sativa</i> (Rapides cultivar-group)      |
| 748   | 3   | AAK74000  | 1337782  | -  | <i>Oryza sativa</i> (Rapides cultivar-group)      |
| 577   | 2   | AAK74546  | 2051606  | -  | <i>Thermomyces blattae</i>                        |
| 144   | 0   | AAK74546  | 2051606  | -  | <i>Thermomyces blattae</i>                        |
| 523   | 2   | BAE6007   | 1017490  | -  | <i>Asciella heloburan</i>                         |
| 538   | 2   | BAB42321  | 1370109  | -  | <i>Staphylococcus aureus</i> subsp. <i>aureus</i> |
| 526   | 2   | BAF95044  | 3120436  | -  | <i>Staphylococcus aureus</i> subsp. <i>aureus</i> |
| 526   | 2   | BAF95044  | 3120436  | -  | <i>Staphylococcus aureus</i> subsp. <i>aureus</i> |
| 526   | 2   | AAK32988  | 1772320  | -  | <i>Pectenovula helicina</i>                       |
| 524   | 2   | AAK32988  | 1772320  | -  | <i>Escherichia coli</i>                           |
| 523   | 2   | AAK32988  | 1772320  | -  | <i>Escherichia coli</i> O157:H7 EDL933            |
| 523   | 2   | AAK32988  | 1772320  | -  | <i>Escherichia coli</i> O157:H7 EDL933            |
| 521   | 2   | AAK11003  | 7326665  | -  | <i>Neisseria meningitidis</i> MC58                |
| 519   | 2   | AAK11003  | 7326665  | -  | <i>Neisseria meningitidis</i> MC58                |
| 523   | 2   | CAB14893  | 7300287  | -  | <i>Neisseria meningitidis</i> 12491               |
| 521   | 2   | AAK78930  | 13024811 | -  | <i>Yersinia pestis</i>                            |
| 517   | 2   | AAK78930  | 13024811 | -  | <i>Clostridium acetobutylicum</i>                 |
| 517   | 2   | CAC89230  | 13978469 | -  | <i>Yersinia pestis</i>                            |
| 512   | 2   | BAE60002  | 15164657 | -  | <i>Clostridium perfringens</i> str. 13            |
| 512   | 2   | BAE60002  | 15164657 | -  | <i>Clostridium perfringens</i> str. 13            |

Finding Modeling Template

NCBI

BLAST Genome Protein Structure PubMed Taxonomy Help

Query: g455752 (NM\_000249) mutL homolog 1; mutL (E. coli) homolog 1; mutL (E. coli) homolog 1 (colon cancer, nonpolyposis type 2) [Homo sapiens]

Matching #453989 631299 730028 741681 1079787 13805126

BlastThis Common Tree Taxonomy Report 3D structures CDD-Search GBL

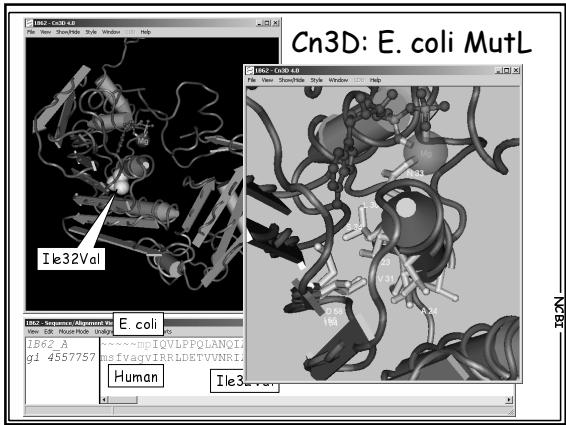
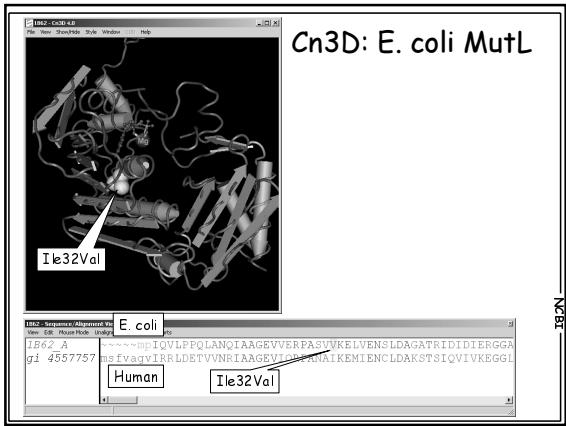
6 BLAST hits to 2 unique species. Sort by taxonomy proximity.

Archaea Bacteria Metazoa Fungi Plants Viruses Other Eukaryotes

Keep only Cut-Off 100 Select Reset

T56 aa

| SCORE | P | ACCESSION | GI       | PROTEIN DESCRIPTION   |
|-------|---|-----------|----------|---|
| 523   | * | 1B62A     | 4929884  | Chain A. MutL Complexed With Adp  |
| 523   | * | 1B62A     | 5542073  | Chain A. MutL Complexed With Adp. N-Terminal 40kDa Fragment Of R          |
| 520   | * | 1B62A     | 5542073  | Chain A. MutL Complexed With Adp. N-Terminal 40kDa Fragment Of Human Pmc2 |
| 421   | * | 1W76A     | 17945768 | Chain A. N-Terminal 40kDa Fragment Of Mcm2 Complexed With Adp             |
| 420   | * | 1W76A     | 17945778 | Chain A. Rgs2-Atppg   |



**Basic Local Alignment Search Tool**

**BLAST**

What's NEW in BLAST®

- March 6th 2002: New database module from BLAST results. Results of a BLAST search will now sequences from the BLAST results page to the NCBI local jackso UnGene databases. Links to additional database soon.

**Nucleotide BLAST**

- Standard nucleotide-nucleotide BLAST (blastn)
- MEGABLAST
- Search for short nearly exact matches

**Protein BLAST**

- Standard protein-protein BLAST (blastp)
- PST and PHL-BLAST
- Search for short nearly exact matches

**Translated BLAST Searches**

- Nucleotide query - Protein db [blastx]
- Protein query - Translated db [tblastx]
- Nucleotide query - Translated db [tblastc]

**Search for conserved domains**

- Search the Conserved Domain Database using RPS-BLAST
- Search by domain architecture [DART]

**Pairwise BLAST**

- BLAST 2 Sequences

**Genomic BLAST pages**

- Human Genome
- Mouse Genome
- Rat Genome
- Fugu rubripes
- Zebrafish Genome
- Anopheles gambiae
- Arabidopsis thaliana
- Oryza sativa
- Other eukaryotes
- Microbial Genomes

Finding Human ESTs

**NCBI** Nucleotide Protein Translations Retrieve results for an RID

Search:

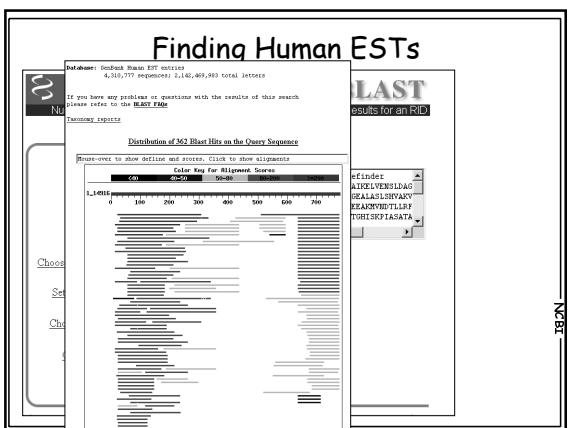
Choose a translation: PROTEIN query - TRANSLATED database [Blastn]

Set subsequence: From  To

Choose database: test\_human

Genetic codes: Disabled

Now: **(BLAST)** or [REVERSE QUERY](#) [REFRESH](#)



Finding Human ESTs

Sequences producing significant alignments:

| Score | E     | bits | Value |
|-------|-------|------|-------|
| 222   | 3e-14 | 22   |       |
| 222   | 5e-58 | 227  |       |
| 221   | 3e-56 | 221  |       |
| 221   | 3e-56 | 221  |       |
| 216   | 2e-55 | 216  |       |
| 181   | 1e-52 | 181  |       |
| 181   | 1e-52 | 181  |       |
| 130   | 2e-50 | 130  |       |
| 201   | 4e-50 | 201  |       |
| 198   | 1e-49 | 198  |       |
| 192   | 6e-48 | 192  |       |
| 131   | 2e-47 | 131  |       |
| 143   | 3e-45 | 143  |       |
| 183   | 6e-45 | 183  |       |
| 180   | 6e-44 | 180  |       |
| 125   | 2e-43 | 125  |       |
| 175   | 9e-42 | 175  |       |
| 175   | 9e-42 | 175  |       |
| 174   | 5e-42 | 174  |       |
| 173   | 8e-42 | 173  |       |
| 167   | 1e-40 | 167  |       |
| 167   | 4e-40 | 167  |       |
| 160   | 4e-40 | 160  |       |
| 164   | 4e-39 | 164  |       |
| 123   | 4e-39 | 123  |       |

Choose a translation:

Set subsequence:

Choose database:

Genetic codes:

Now: **(BLAST)** or [REVERSE QUERY](#) [REFRESH](#)

**UniGene Cluster HG57801 Homo sapiens**

**MLH1 MutL homolog 1, colon cancer, nonpolyposis type 2 (E. coli)**

**SEE ALSO**

- GeneLink: 4262
- OMIM: 16
- HomoloGene: 1

**SELECTED MUTL PROTEIN SIMILARITIES**

| Organism, protein | Percent identity and length of aligned protein | Link                         |
|-------------------|--|------------------------------|
| (Human)           | ->0000001 ->MTL PROTEIN HOMOLOG 1 (HOMOL0001)  | <a href="#">View Details</a> |

**KnownGenes:**

- g27072 - MLH1\_MOUSE DNA mismatch repair protein* 96% / 754 aa (see [Protein](#))
- pr715102 - TS1520 DNA mismatch repair protein* 99% / 748 aa (see [Protein](#))
- Achilles*: *pr715102 - TS1520 DNA mismatch repair protein* 99% / 748 aa (see [Protein](#))
- CatGenes*: *pr715102 - TS1520 DNA mismatch repair protein* 99% / 748 aa (see [Protein](#))
- Ecoli*: *pr715102 - TS1520 DNA mismatch repair protein* 99% / 748 aa (see [Protein](#))
- Saccharomyces cerevisiae*: *MLH1\_YEAST MUTL PROTEIN* 97% / 748 aa (see [Protein](#))

**MAPPING INFORMATION**

Chromosome: 3  
Genomic View: [MutL homolog 1](#)  
Genomic Gene Map: [3.2](#)

Whitehead map: [WJ7345](#), Chr 3, YAC contig WJ7345  
Units entry: [SGD.GC.L2125](#) Genetics Context Map View

**EXPRESSION INFORMATION**

DNA sources: 2 pooled whole human, one primary and one metastatic to breast, bone, brain, liver, lung, melanoma, skin, stomach, testis, uterus (lymphoma/blood/bone/brain/testis/liver/normal/breast/lymphoma/cervix, cell line/choriocarcinoma/colon)

---

**UniGene Cluster HG57801 Homo sapiens**

**MLH1 MutL homolog 1, colon cancer, nonpolyposis type 2 (E. coli)**

**SEE ALSO**

- GeneLink: 4262
- OMIM: 16
- HomoloGene: 1

**SELECTED MUTL PROTEIN SIMILARITIES**

| Organism, protein | Percent identity and length of aligned protein | Link                         |
|-------------------|--|------------------------------|
| (Human)           | ->0000001 ->MTL PROTEIN HOMOLOG 1 (HOMOL0001)  | <a href="#">View Details</a> |

**KnownGenes:**

- g27072 - MLH1\_MOUSE DNA mismatch repair protein* 96% / 754 aa (see [Protein](#))
- MLH1 (Mutl homolog 1)*: *pr715102 - TS1520 DNA mismatch repair protein* 99% / 748 aa (see [Protein](#))
- Achilles*: *pr715102 - TS1520 DNA mismatch repair protein* 99% / 748 aa (see [Protein](#))
- CatGenes*: *pr715102 - TS1520 DNA mismatch repair protein* 99% / 748 aa (see [Protein](#))
- Ecoli*: *pr715102 - TS1520 DNA mismatch repair protein* 99% / 748 aa (see [Protein](#))
- Saccharomyces cerevisiae*: *MLH1\_YEAST MUTL PROTEIN* 97% / 748 aa (see [Protein](#))
- Homologs*: *g27072 - MLH1\_MOUSE DNA mismatch repair protein* 96% / 754 aa (see [Protein](#))

**MAPPING INFORMATION**

Chromosome: 3  
Genomic View: [Mutl homolog 1](#)  
Genomic Gene Map: [3p13](#)

Whitehead map: [WJ7345](#), Chr 3, YAC contig WJ7345  
Units entry: [SGD.GC.L2125](#) Genetics Context Map View

**EST SEQUENCES (10 of 278) [Show all ESTs]**

| EST Sequence | Description   | Score | Read Type |
|--------------|---|-------|-----------|
| BE888481     | CDNA clone IMAGE:9911872 (leiomysarcoma)                  | 99    | read      |
| BE327277     | CDNA clone IMAGE:4123836 (retal cell adenocarcinoma)      | 99    | read      |
| BE306509     | CDNA clone IMAGE:4123836 (rhabdomyosarcoma)               | 99    | read      |
| BE264482     | CDNA clone IMAGE:4123806 (choriocarcinoma)                | 99    | read      |
| BE795996     | CDNA clone IMAGE:4343565 (lymphoma, cell line)            | 99    | read      |
| BE850881     | CDNA clone IMAGE:5226825 (pooled pancreas and spleen)     | 99    | read      |
| BE539316     | CDNA clone IMAGE:3451538 (placenta)                       | 99    | read      |
| BE841102     | CDNA clone IMAGE:0013904 (epidermod carcinoma, cell line) | 99    | read      |
| BE9772547    | CDNA clone IMAGE:4837570 (testis, cell line)              | 99    | read      |
| BE9772733    | CDNA clone IMAGE:4837559 (testes, cell line)              | 99    | read      |

---

**LocusLink MLH1**

**Links:**

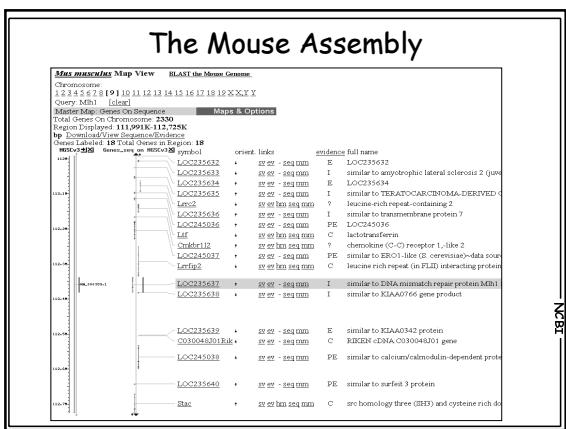
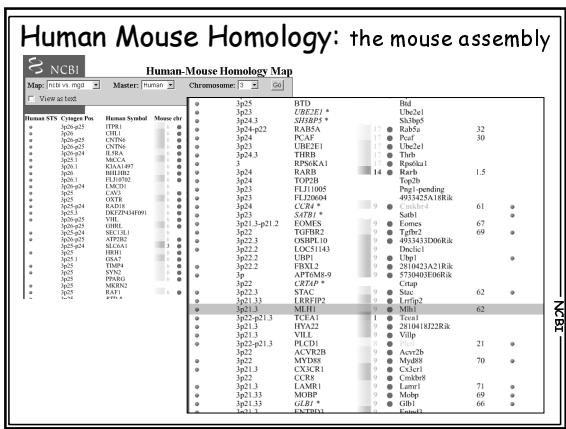
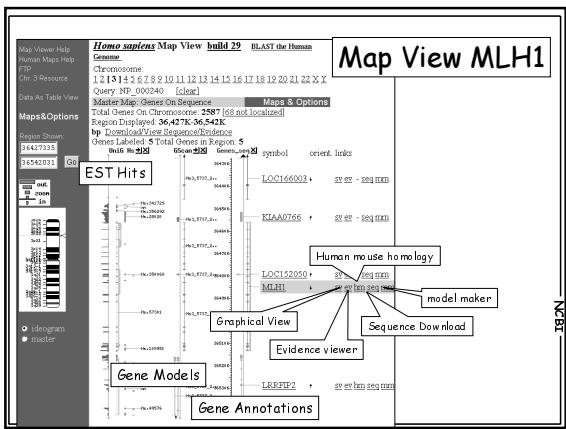
| pm                   | PubMed               |
|----------------------|----------------------|
| <a href="#">View</a> | <a href="#">View</a> |
| mv                   | MapViewer            |
| <a href="#">View</a> | <a href="#">View</a> |
| sv                   | Sequence Viewer      |
| <a href="#">View</a> | <a href="#">View</a> |
| ev                   | Evidence Viewer      |
| <a href="#">View</a> | <a href="#">View</a> |
| BL                   | BLASTLink            |
| <a href="#">View</a> | <a href="#">View</a> |

**NCBI Reference Sequence (RefSeq)**

Category: REVIEWED  
mRNA: NM\_000249  
Name: MLH1, mutl homolog 1  
Description: Human: Human homolog of the E. coli mutL homolog 1. Human: Human homolog of the E. coli mutL homolog 1. Human: Human homolog of the E. coli mutL homolog 1. Human: Human homolog of the E. coli mutL homolog 1. Human: Human homolog of the E. coli mutL homolog 1. Human: Human homolog of the E. coli mutL homolog 1. Human: Human homolog of the E. coli mutL homolog 1. Human: Human homolog of the E. coli mutL homolog 1. Human: Human homolog of the E. coli mutL homolog 1.

Category: NCBI Genome Annotation  
Genomic Contig: NT\_027530  
Annotated transcripts present for this locus:  
Evidence: supported by alignment with mRNA  
mRNA: NM\_000249  
mRNA: NT\_027530

---



## Live Web Demos

Genetic mapping Fanconi renal syndrome

Finding Prestin Homologues

NB:

---

---

---

---

---

---

---